

CONSIDERATIONS ON THE GENETIC EQUILIBRIUM LAW

SIMONE CAMOSSO

ABSTRACT. In the first part of the paper I will present a brief review on the Hardy–Weinberg equilibrium and its formulation in projective algebraic geometry. In the second and last part I will discuss examples and generalizations on the topic.

CONTENTS

1. Introduction	1
2. The Hardy–Weinberg law	1
3. Projective and algebraic geometry	2
4. Examples, generalizations and conclusion	3
References	4

1. INTRODUCTION

The study of population genetics, evolution and its evolutionary trees are classical subjects in biology. A mathematical approach consists in the maximum likelihood estimation, a technique largely used in statistic. This approach leads to the problem of maximizing particular functions of certain parameters. A theoretical study and an upper bound for the maximum likelihood degree is discussed in [3]. Different techniques in order to solve likelihood equations are described in [9]. Applications of these methods have been used when the statistical model is an algebraic variety, this is the case of Fermat hypersurfaces treated in [1]. In the other side we have applications of these ideas to biological models. In this direction we refer to a work of [4] where phylogenetic models in two different topologies have been studied by the authors.

The purpose of this paper is a “soft” introduction to these ideas with a discussion on the Hardy–Weinberg case.

2. THE HARDY–WEINBERG LAW

The Hardy–Weinberg law states that allele and genotype frequencies in a population remain constant during the generation change. This happens under the following assumptions: the size of the population must be very large, we have absence of migration and mutations, the mating is random and the natural selection doesn’t affect the alleles under consideration. Mathematically if p represents the number of pure dominants characters AA , q the number of heterozygotes Aa and r the number of pure recessives aa , the following proportion holds $p : 2q : r$ (see

[7]). Another way, if p and q represent the allele frequencies of the character A and a with $p + q = 1$, taking the square we find that:

$$p^2 + 2pq + q^2 = 1, \quad (2.1)$$

where p^2 , $2pq$ and q^2 represent the genotype frequencies associated to AA , Aa and aa . The equation (2.1) describes the constancy of the genotypic composition of the population and is called the Hardy–Weinberg principle or the Hardy–Weinberg equilibrium (HWE). We consider [6] and [8] as scholarly references on this subject. Different generalizations of (2.1) are possible. The first concern the number of alleles at a locus. For example in the case of three alleles A_1, A_2 and A_3 , with frequencies respectively given by p, q and r , the genotype frequencies are given by the following expansion $(p + q + r)^2 = p^2 + q^2 + r^2 + 2pq + 2qr + 2pr$. In general, for any number n of alleles with frequencies p_i , we have that:

$$(p_1 + \dots + p_n)^2 = \sum_{i=1}^n p_i^2 + \sum_{i \neq j} 2p_i p_j.$$

In another direction the generalization is given considering the binomial $(p + q)^m$, with $m = 3, 4, 5, \dots$. This is the case of polyploid. For example considering tetraploids ($m = 4$) the procedure involves the expansion of $(p + q)^4$. We observe how in this particular example the frequency of heterozygotes (given by the mixed terms in the expansion) is $2pq(2 - pq)$ that is considerably greater than $2pq$, the frequency for a diploid organism. More information on this topic can be found in [6].

3. PROJECTIVE AND ALGEBRAIC GEOMETRY

In this section we shall examine how translate previous considerations in the modern language of projective and algebraic geometry. The setting is the same of [11] in its first lecture. We shall show how the Hardy–Weinberg law can be formulated in a fixed system of homogeneous coordinates l, m, n in \mathbb{P}^2 . First we shall consider the open triangle $\Delta_2 = \{(l, m, n) \in \mathbb{R}_+^3 : l + m + n = 1\}$, where \mathbb{R}_+ are the positive reals. Second we shall observe that setting $l = p^2, m = 2pq, n = q^2$ to be genotype frequencies, we have the relation:

$$m^2 = 4ln, \quad (3.1)$$

that is the equation of a parabola in the triangle of vertex $(1, 0, 0), (0, 1, 0), (0, 0, 1)$. We call the zero locus of (3.1), denoted also by $V(m^2 - 4ln)$, the Hardy–Weinberg curve (details are in [5]). In the theory of [11] (and [12]) is of particular interest a function called “likelihood function” that depends by some positive integer parameters. This function is positive on Δ_2 and zero on the boundary of Δ_2 . We shall denote this function by l and with u_0, u_1, \dots the corresponding parameters. In the case of the Hardy–Weinberg curve this function has the following form:

$$l_{u_0, u_1, u_2} = l^{u_0} m^{u_1} n^{u_2} = 2^{u_1} p^{2u_0 + u_1} q^{u_1 + 2u_2}.$$

We observe that l_{u_0, u_1, u_2} is a function depending only by the variable p (because $q = 1 - p$) and the MLE problem consists in the estimation of p maximizing the function l_{u_0, u_1, u_2} . Lagrange Multipliers can be used in order to solve the problem and in this case the solution is given by the point:

$$\hat{p} = \frac{2u_0 + u_1}{2u_0 + 2u_1 + 2u_2}. \quad (3.2)$$

4. EXAMPLES, GENERALIZATIONS AND CONCLUSION

As exercise we shall apply the same procedure in order to solve the MLE problem for the case of three alleles and in the second time for the case of tetraploids. For the first we shall start writing the “likelihood function” associated to the HWE given by $(p + q + r)^2$, where p, q and r are the usual frequencies. In this case we have that:

$$l_{u_0, u_1, u_2, u_3, u_4, u_5} = 2^{u_3+u_4+u_5} p^{2u_0+u_3+u_5} q^{2u_1+u_3+u_4} (1-p-q)^{2u_2+u_4+u_5}.$$

We shall proceed maximizing the function of two variables p, q . This is an ordinary problem of calculus that gives as answer the point:

$$\left(\frac{2u_0 + u_3 + u_5}{2u_1 + 2u_2 + 2u_4 + u_3 + u_5}, \frac{u_4 - u_5 + 2u_1 - 2u_0}{2u_1 + 2u_2 + 2u_4 + u_3 + u_5}, \frac{2u_2 + u_4 + u_5}{2u_1 + 2u_2 + 2u_4 + u_3 + u_5} \right).$$

For the tetraploid case, before to proceed, we shall observe that calling $l_0 = p^4, l_1 = q^4, l_2 = 4pq^3, l_3 = 4p^3q, l_4 = 6p^2q^2$ the genotype frequencies, the Hardy–Weinberg equilibrium can be represented by the following relation:

$$l_4^4 = \frac{1}{81} l_0 l_1 l_2 l_3.$$

The associated “likelihood function” is

$$l_{u_0, u_1, u_2, u_3, u_4} = l_0^{u_0} l_1^{u_1} l_2^{u_2} l_3^{u_3} l_4^{u_4}.$$

We shall make the expedient of consider the logarithm of the previous function instead the original finding as maximizing point:

$$\left(\frac{u_0}{|u|}, \frac{u_1}{|u|}, \frac{u_2}{|u|}, \frac{u_3}{|u|}, \frac{u_4}{|u|} \right),$$

where $|u| = u_0 + u_1 + u_2 + u_3 + u_4$. We recommended the use of a scientific software, as Maple or MATLAB, especially when the number of parameters is considerably high.

Now I want to spend these last words comparing analogies between the HWE and the algebraic geometry. It is clear that a possible extension of the Hardy–Weinberg law can take the following form:

$$(p_0 + \dots + p_n)^m = c, \quad (4.1)$$

where c is some constant and p_i from $i = 0, \dots, n$ are the allele frequencies such that the sum is fixed. Now expanding (4.1) we find the polynomial form:

$$\sum_{i_0 + \dots + i_n = m} \frac{m!}{i_0! \dots i_n!} p_0^{i_0} \dots p_n^{i_n} = c.$$

From the algebraic geometry point of view this is the image of the Veronese map $\nu_m : \mathbb{P}^n \rightarrow \mathbb{P}^{\binom{n+m}{m}-1}$ given by $(p_0, \dots, p_n) \mapsto (p_0^m, p_0^{m-1}p_1, \dots, p_n^m)$ (see [2]). The classical Hardy–Weinberg law corresponds to the case of $\nu_1 : \mathbb{P}^1 \rightarrow \mathbb{P}^2$ that

$(p, q) \mapsto (p^2, pq, q^2)$ and $c = 1$. Using the identification between homogeneous polynomials that are power of linear forms and the image of the Veronese map, we can think these generalized laws as Veronese projective varieties. From the side of algebraic geometry there are a rich collection of results concerning the Veronese and Segre varieties, for example it is possible to compute the Hilbert polynomial and other invariants. It is not all peace and light because the constraint $\Delta_n = \{(p_0, \dots, p_n) \in \mathbb{R}_+^{n+1} : p_0 + \dots + p_n = 1\}$ doesn't permit the complete translation of the problem using the previous identification.

Anyway the methods of numerical algebraic geometry seem to give good prospects in this direction and in [10] the ML degree has been calculated for matrices with rank constraints. In particular the case of rank one gives the ML degree equal to one, so \hat{p} is a rational function of a set of parameters u_0, u_1, \dots .

REFERENCES

- [1] D.Agostini, D.Alberelli, F.Grande, P.Lella, "The maximum likelihood degree of Fermat hypersurfaces", arXiv:1404.5745.
- [2] E.Arrondo, "Introduction to projective varieties", unpublished notes from the website: <http://www.mat.ucm.es/~arrondo/projvar.pdf> (2007), 9–10.
- [3] F.Catanese, S.Hoten, A.Khetan, B.Sturmfels, "The maximum likelihood degree", Amer. J. Math. 128 (2006), no. 3, 671–697. MR 2230921.
- [4] B.Chor, A.Khetan, S.Snir, "Maximum likelihood on four taxa phylogenetic trees: analytic solutions", The 7th Annual Conference on Research in Computational Molecular Biology–RECOMB 2003, Berlin, April 2003, pp. 76–83.
- [5] A.W.F.Edwards, "Foundations of Mathematical Genetics", Cambridge University Press, Cambridge (2000).
- [6] R.Frankham, J.D.Ballou, D.A.Briscoe, "Introduction to conservation genetics", Cambridge (2002), 86–90.
- [7] G.H.Hardy, "Mendelian Proportions in a Mixed Population", Science, New Series, Vol.28, 706 (1908), 49–50.
- [8] D.L.Hartl, A.G.Clark, "Principles of population genetics", Sinauer Associates, Inc. Publishers, Sunderland, Massachusetts (Fourth Edition).
- [9] S.Hoşten, A.Khetan, B.Sturmfels, "Solving the Likelihood Equations", Foundations of Computational Mathematics, Vol. 5, Issue 4, pp 389–407, 2005.
- [10] J.Hauenstein, J.Rodriguez, B.Sturmfels, "Maximum Likelihood for Matrices with Rank Constraints", Journal of Algebraic Statistics, Vol.5, Issue 1 (2014), pp 18–38.
- [11] J.Huh, B.Sturmfels, "Likelihood Geometry", Combinatorial Algebraic Geometry: Levico Terme, Italy 2013, Springer International Publishing, Vol.2108 of the series Lecture Notes in Mathematics (2014), 63–117.
- [12] I.J.Myung, "Tutorial on maximum likelihood estimation", Journal of Mathematical Psychology 47 (2003) 90–100.